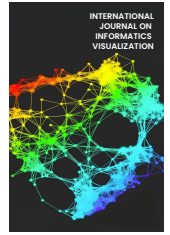




INTERNATIONAL JOURNAL ON INFORMATICS VISUALIZATION

journal homepage : www.joiv.org/index.php/joiv



Classification of Predicting Customer Ad Clicks Using Logistic Regression and k -Nearest Neighbors

Yasi Dani ^{a,*}, Maria Artanta Ginting ^a

^a Computer Science Department, School of Computer Science, Bina Nusantara University, Bandung Campus, Jakarta 11480, Indonesia

Corresponding author: *yasi.dani@binus.ac.id

Abstract—Nowadays, conventional marketing techniques have changed to online (digital) marketing techniques requiring internet access. Online marketing techniques have many advantages, especially in terms of cost efficiency and fast information delivery to the public. Therefore, many companies are interested in online marketing and advertising on social media platforms and websites. However, one of the challenges for companies in online marketing is determining the right target consumers since if they target consumers who are not interested in buying the product, the advertising costs will be high. One use of online advertising is clicks on ads which is a marketing measurement of how many users click on the online ad. Thus, companies need a click prediction system to know the right target consumers. And different types of advertisers and search engines rely on modeling to predict ad clicks accurately. This paper constructs the customer ad clicks prediction model using the machine learning approach that becomes more sophisticated in effectively predicting the probability of a click. We propose two classification algorithms: the logistic regression (LR) classifier, which produces probabilistic outputs, and the k -nearest neighbors (k -NN) classifier, which produces non-probabilistic outputs. Furthermore, this study compares the two classification algorithms and determines the best algorithm based on their performance. We calculate the confusion matrix and several metrics: precision, recall, accuracy, F1-score, and AUC-ROC. The experiments show that the logistic regression algorithm performs best on a given dataset.

Keywords—Machine learning algorithm; logistic regression; k -nearest neighbors; supervised classification; ad clicks.

Manuscript received 6 Jul. 2022; revised 30 Oct. 2022; accepted 16 Nov. 2022. Date of publication 31 Mar. 2023.
International Journal on Informatics Visualization is licensed under a Creative Commons Attribution-Share Alike 4.0 International License.



I. INTRODUCTION

The Internet is a communication network connecting computer networks with other devices for users worldwide. The Internet is one proof of technological sophistication at this time, where this technology is believed to be the sixth most powerful medium of all existing media. The impact of the Internet is positive or negative; it depends on the user using it. In general, the positive effects are to support and facilitate all activities of modern society since all activities can be done online. This is the main trigger in advertising to switch to online advertising. Previously, advertising media was done traditionally, such as magazines, brochures, newspapers, television, radio, and so forth [1]. The main goal of online advertisement is its ability to reach the world community without geographical boundaries with a wide reach easily and to save marketing costs. This online advertisement is a new strategy to increase website traffic and deliver advertisements to the right customers [2].

In recent years, people have been affected by the biggest coronavirus disease 2019 (COVID-19) outbreak worldwide. As a result of this pandemic, online consumers have increased significantly due to physical distancing rules and staying at home so that conventional consumers are forced to shop online. Hence, conventional shopping is very ineffective at this time [3]. Ultimately, consumers spend more time and money than ever before on online sites [4]. According to a report in the United States, online advertising is increasing rapidly where online advertising revenue increased from 2019 to 2020 which increased by about 12.2% (from \$124.6 billion grew to \$139.8 billion), then online advertising in 2025 is predicted will reach \$982.82 billion, so that according to statistical data in the United States that online advertising revenue is greater than conventional advertising. Online customers have also increased sharply since online advertisements promote their products using websites to deliver attractive advertisements to consumers [5]. Eventually, online advertising spreads in various countries because of the wide internet network and many online advertisers. Thus, online advertising in many countries has become a new trend

in marketing and business fields because of its effectiveness and efficiency [6]. However, it could not be denied that online advertising also has big challenges today, especially in the demands of creativity and innovation. The new focus on online advertising trends that must be different and unique. This is a big challenge for entrepreneurs in creating online applications and developing social sites [7]. Therefore, online advertising has grown rapidly in recent years. This phenomenon becomes a great opportunity for future research fields [8].

One of the weaknesses in conventional advertising methods is that advertising methods are often not right on target and this method does not directly lead to the right target market and consumers. Whereas in the online advertising method, advertisers can allocate their ads dynamically, that is, they can adjust directly to user interests based on the feedback they observe. This online advertising business is very promising due to the increasing income every year [9]. Online advertising has proven to be cheaper than conventional advertising. And online advertising can now be directed to the desired target consumers. Therefore, studies are needed to develop methods to predict ad clicks. Ad click prediction has been used historically in every type of ad format, like textual and contextual ads, search engines, video ads, etc. Currently, advertising is increasing rapidly, and ad click prediction requires a lot of data analysis [10].

Machine learning is one area of artificial intelligence (AI) that has created many things. This learning is based on mathematics and statistics, which aim to assist in future predictions by analyzing historical data [11]. Currently, many machine learning algorithms have been used by many researchers, one of them being predictive analysis. Machine learning is important in calculating the expected utility of potential advertisements to customers and the efficiency of marketing costs. This advertising efficiency depends on the level of accuracy in predicting ad clicks, where the prediction analysis must be robust and adaptive in the prediction model since the volume of data will be very large [12]. Furthermore, the chosen methodology will bring changes in the prediction analysis result of ad clicks [13].

One of the domains of machine learning is predictive analysis, advertising prediction is considered to be one of the most profitable things. Machine learning is an ideal learning method to assist management in making decisions based on data where system and logic in machine learning are highly recommended by statisticians because they are not only based on historical data but also consider other parameters. Furthermore, in machine learning the chosen methodology and technique will bring about changes in the prediction analysis of ad clicks.

In previous research, several machine learning techniques have been applied to predict clicks on advertisements such as Richardson et al. [14] proposed a logistic regression classifier to estimate the click-through rate for new ads, Cheng et al. [15] proposed maximum entropy in sponsored search, Broder et al. [16] studied support vector machines (SVM) to predict whether or not to show any of the ads for the incoming request, Guo et al. [17] proposed a conditional random field (CRF) to detect the user's search goals, and Chakrabarti et al. [18] applied logistic regression to learn the match between ads and Web pages.

In this research, we propose a model using two supervised learning algorithms, the logistic regression, and the k-Nearest Neighbors classifier, to predict customer ad clicks. This paper aims to help companies use the right method for predictive analysis to reach the right target consumers. Our algorithm performances are evaluated using some metrics that are precision, recall, accuracy, F1-score, and AUC-ROC.

The remainder of the paper is organized as follows: Section two explains the methods of the logistic regression classifier and *k*-nearest neighbors' classifier and how to evaluate the performance of both algorithms on a given dataset. In section three, we give the results of the performance of the two algorithms on the dataset, then compare the performance and determine the best algorithm, and a conclusion is presented in Section 4.

II. MATERIALS AND METHOD

In this section, we present some classification algorithms in machine learning and the methodology of this research.

A. Logistic Regression (LR)

Logistic regression (LR) is a method that is often used for classification, which is a statistical analysis technique applied for predictive models [19], [20]. This classification is one of the most popular machine learning algorithms that come under supervised learning techniques. Moreover, this classification model usually achieves high algorithm performance, so it is often applied in the industrial world [21], [22]. There are several types of logistic regression, namely binary and multinomial logistic regression [23], [24]. Binary logistic regression is used when the response variable is dichotomous. That is, there are only two categories [25]. Meanwhile, multinomial linear regression is used when the response variable has more than two categories [26]. This research uses binary linear regression.

Another advantage of the logistic regression model is the ability to process large volumes of data at high speed because it requires less computational capacity, such as memory and processing power. This makes the model very suitable for data scientists to get multiple solutions with fast results. Logistic regression is also used extensively in the fields of medicine and social sciences, as well as in marketing, such as predicting a customer's propensity to buy a product or unsubscribe.

This logistic regression is a predictive model similar to linear regression based on the logistic function or the Sigmoid function [27]. The difference between the results of linear regression and logistic regression is that the range of values in linear regression is a real number, while the range of values in logistic regression is between 0 and 1. Then it also does not require a linear relationship between input and output variables since it uses a nonlinear log transformation approach to predict the odds ratio [28], [29]. In general, the assumptions of LR include:

1. There is no need for linearity between the independent and response variables.
2. There is no need to assume multivariate normality or equal variance between independent variables.
3. There is no need for the assumption of homoscedasticity.
4. The dependent variable must be dichotomous.
5. Do not need to transform into metric form.

6. The categories must be separate or exclusive to the independent variables.
7. Requires a relatively large sample for predictor variables, for example, a minimum of 50 data samples.
8. The odds ratio is a probability value.

Given vector data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ where n is the number of instances and d is the number of features (parameters), and y be a binary outcomes vector. And the vector θ is the vector of unknown parameters such that $\mathbf{x}_i \leftarrow [1, \mathbf{x}_i]$ and $\theta \leftarrow [\theta_0, \theta^T]$. From now on, the assumption is that the intercept is included in the vector θ . For every instance $\mathbf{x}_i \in \mathbb{R}^d$ where $i = 1, 2, \dots, n$, the outcome is either $y_i = 1$ (positive instance) or $y_i = -1$ (negative instance). The logistic function commonly used to model each positive instance \mathbf{x}_i with its expected binary outcome is given by

$$E[y_i = 1 | \mathbf{x}_i, \theta] = p_i = \frac{e^{\mathbf{x}_i \theta}}{1 + e^{\mathbf{x}_i \theta}} = \frac{1}{1 + e^{-\mathbf{x}_i \theta}}, \quad (1)$$

where $i = 1, 2, 3, \dots, n$.

B. K-Nearest Neighbors (k-NN)

Nearest neighbor search is one of the most popular learning in the field of machine learning and the classification technique introduced by Fix and Hodges. This learning has proven to be a simple and powerful recognition algorithm. Cover and Hart show that decision rules work well given that no explicit knowledge about the data is available. A simple generalization of this method is called the k -NN rule, in which new patterns are classified into the class with the most members among the k -nearest neighbors. This can be used to obtain a good estimate of Bayes error and the probability of error is asymptotically close to Bayes error.

In the classification technique, the different characteristics in the classification determine the class where the unlabeled data resides with the aim of classifying data based on the closest or neighboring training examples in a particular region. The advantage of this technique is the simplicity of execution and low computation time. For continuous data, it uses the Euclidean distance to calculate its nearest neighbors.

K -nearest neighbor (k -NN) is one of the statistical analysis techniques to build the simplest prediction model since it does not require mathematical assumptions and heavy machinery [30]. K -NN is a non-parametric supervised learning method and is commonly used for classification [31]. K -NN is very popular due to its simplicity and excellent empirical performance. It can handle both binary data and makes no assumptions about the parametric of the decision boundary [32], [33]. This classifier aims to predict the target/class of an observation point based on the closest neighbor k class [34]. In calculating the k -NN method, it takes several steps; namely, we choose the value of k , then calculate the distance with the distance function from one observation to all other observations and take k nearest neighbors as per the calculated distance. After that count the number of observation points in each category among k this neighbor. Ultimately, we assign the new observation point to the category with the most neighbors [35] [36].

Given vector data $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ where n is the number of instances and d is the number of features (parameters), and y be a binary outcomes vector. The goal in classification is to learn a functional model f that allows a

reasonable prediction of class label y' for an unknown pattern \mathbf{x}' . K -NN assigns the class label of the majority of the k -nearest patterns in data space. For this sake, we have to be able to define a similarity measure in data space. In \mathbb{R}^d , it is reasonable to employ the Minkowski metric (p -norm)

$$\|\mathbf{x}' - \mathbf{x}_j\|^p = \left(\sum_{i=1}^d |(x'_i) - (x_i)_j|^p \right)^{\frac{1}{p}}, \quad (2)$$

which corresponds to the Euclidean distance for $p = 2$. In other data spaces, adequate distance functions have to be chosen, e.g., the Hamming distance in \mathcal{B}^d . In the case of binary classification, the label set $Y = \{-1, 1\}$ is employed, and k -NN is defined as

$$f(\mathbf{x}') = \begin{cases} 1 & \text{if } \sum_{i \in M_k(\mathbf{x}')} y_i \geq 0 \\ -1 & \text{if } \sum_{i \in M_k(\mathbf{x}')} y_i < 0 \end{cases} \quad (3)$$

with neighborhood size k and with the set of indices $M_k(\mathbf{x}')$ of the k -nearest patterns.

C. Performance Evaluation

In the field of machine learning and computing, evaluating the performance of a classification algorithm is very important. The goal is to measure the performance of an algorithm so that we can consider it in selecting the best algorithm [40]. The input data is grouped into one of two classes in binary classification, which is the simplest and most widely used form. Measuring the performance of the classification model creates a confusion matrix. The output of this confusion matrix can be two or more classes, this research performs a binary classification, so the confusion matrix results are two classes. The confusion matrix aims to compare the classification results of an algorithm with the ground-truth classification results [41].

The representation of the confusion matrix is a matrix table with four combinations of predicted values, and the actual value where the table can be seen in Table I. Suppose there are two classification results, namely positive (labeled 1) and negative (labeled 0), then the four combinations include 1). True Positive (TP) is the amount of positive data that is predicted to be true as positive, 2). False Negative (FN), which is the amount of data that is positive but is predicted to be negative 3). True Negative (TN) where the number of data that is negative and is predicted to be true as negative, and 4). False Negative (FN) is the amount of data that is positive but is predicted to be negative. Next, when a prediction result is a real number, a threshold value of t is needed to distinguish positive and negative classes, after which the confusion matrix can be made [42].

TABLE I
CONFUSION MATRIX

		Predicted Values	
		Positive	Negative
Actual Values	Positive	TP	FN
	Negative	FP	TN

Furthermore, we use the results of the confusion matrix table to evaluate the performance of the machine learning

algorithm for making predictions, namely by calculating the values of precision, recall/sensitivity, specificity, accuracy, and F1-score. For measuring algorithm performance, we could calculate some metrics that are sensitivity or recall (Eq. 3), specificity (Eq. 4), precision (Eq. 5), accuracy (Eq. 6), and F1-score (Eq.7) [43] :

$$\text{sensitivity/recall (true positive rate)} = \frac{TP}{TP+FN} \quad (3)$$

$$\text{specificity (true negative rate)} = \frac{TN}{TN+FP} \quad (4)$$

$$\text{precision (positive predictive value)} = \frac{TP}{TP+FP} \quad (5)$$

$$\text{accuracy} = \frac{TP+TN}{TP+FN+TN+FP} \times 100\% \quad (6)$$

$$\text{F1-score} = \frac{2 \times \text{precision} \times \text{recall}}{\text{precision} + \text{recall}} \quad (7)$$

To add a measure to evaluate the performance of the algorithm, we also use ROC (Receiver Operating Characteristic) analysis which calculates a confusion matrix for all possible threshold values. This ROC curve represents the relationship between the true positive rate or sensitivity (y-axis) and the false-positive rate or 1-specificity (x-axis). After calculating the ROC value, we plot the curve of all ROC values, then calculate the area under the curve, called AUC (Area Under Curve). This AUC-ROC is an area that describes the level of accuracy of the algorithm. The range of AUC-ROC values is between 0 and 1. Generally, the higher the AUC-ROC score, the better a classifier performs for the given task.

D. Dataset

In this research, we use a dataset taken from Kaggle's website about the online advertising of a marketing agency. Then, we process this dataset to predict whether a particular online user will click on an online ad. Therefore, we apply several classification algorithms to predict it. This dataset consists of 1000 observations and 10 features which are: 1). Daily Time Spent on Site, 2). Age, 3). Area Income, 4). Daily Internet Usage, 5). Ad Topic Line, 6). City, 7). Male, 8). Country, 9). Timestamp and 10). Clicked on Ad. The response feature is Clicked on Ad. This feature has two possible outcomes that are 0 and 1 where 0 refers to the case where a user didn't click the advertisement (class 0), while class 1 refers to the scenario where a user clicks the advertisement (class 1). We use features: 'Daily Time Spent on Site' until 'Timestamp' to accurately predict the value 'Clicked on Ad' feature. This research divides data into 67% in training data and 33% in testing data.

E. Methods

There are various methods or classification algorithms in machine learning. Nevertheless, all methods do not have the same accuracy and each algorithm has different accuracy. This research implements two machine learning classification methods on the dataset for predictive analysis and evaluates the performance of each classification algorithm. We use two classification methods or algorithms: logistic regression classifier (LR) and *k*-Nearest Neighbors (*k*-NN) classifier. These methods are very popular in supervised machine learning, so many researchers have good experience with

them since they usually have good algorithm performance. Each of these two algorithms has different steps from each other in classifying. A solution has been developed for the classification of numerical data by these two algorithms.

An architecture overview is shown in Fig. 1. At the beginning, we input the dataset we need to classify. In classifying this data set, we used two machine learning classifiers: LR and *k*-NN. These classifiers were applied to predict if a particular user would click on an online advertisement. Furthermore, to obtain the best classifier method, we evaluate the performance of both classification methods with confusion matrix and several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC.

In the beginning, we input the dataset where the data has been divided into two parts: training data and test data. In classifying this data set, we used two machine learning classifiers: LR and *k*-NN. These classifiers were applied to predict if a particular user would click on the advertisement. Furthermore, to obtain the best classifier method, we evaluate the performance of both classification methods with several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC.

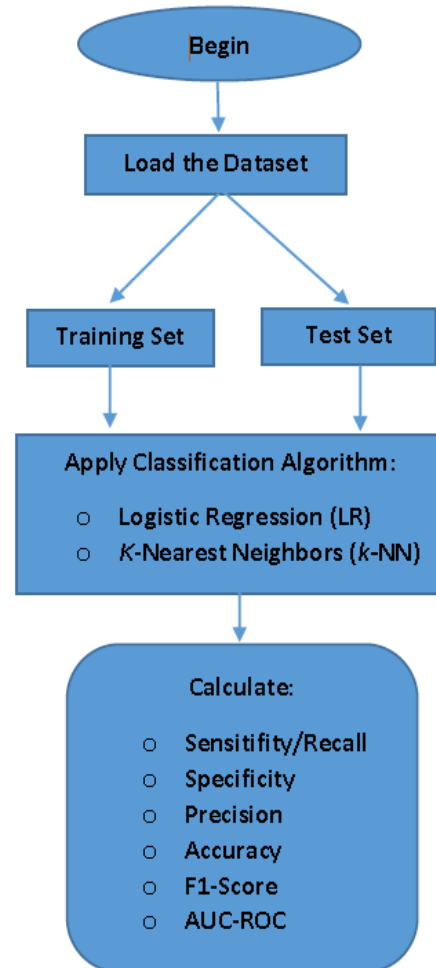


Fig. 1 The Proposed Method Flowchart

III. RESULTS AND DISCUSSION

In this section, we first perform and compare the individual logistic regression and k -NN classification algorithms on the advertisement dataset. And then, we choose the best classification algorithm of the two classification algorithms based on the calculation of several metrics. We use Python 3.7.3 to perform the simulation results of this research.

First, we apply a logistic regression classifier to the training dataset. Then, we evaluate the performance of both classification methods with a confusion matrix and several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC. And then, we apply a logistic regression classifier to the testing dataset. Then, we also evaluate the performance of both classification methods with a confusion matrix and several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC. Next, we apply k -NN algorithm to the training dataset where $k=2$. Then, we also evaluate the performance of both classification methods with a confusion matrix and several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC. And then, we apply k -NN algorithm to the testing dataset where $k=2$. Then, we also evaluate the performance of both classification methods with a confusion matrix and several metrics: sensitivity, specificity, precision, accuracy, F1-score, and AUC-ROC.

TABLE II
COMPARISON OF CONFUSION MATRIX RESULTS FOR TRAINING SET

	Logistic Regression		k -NN Classifier	
	Predicted Class		Predicted Class	
Actual Class	1	0	1	0
1	312	26	274	64
0	43	289	0	332

The tables are the results of two classification methods or algorithms on online advertising datasets. Table II shows that the logistic regression classifier of training set correctly classifies a total of 312 in class 1 and a total of 289 in class 0. And the k -NN algorithm of the training set correctly classifies a total of 274 in class 1 and correctly classifies a total of 332 in class 0.

TABLE III
COMPARISON OF CONFUSION MATRIX RESULTS FOR TESTING SET

	Logistic Regression		K -NN Classifier	
	Predicted Class		Predicted Class	
Actual Class	1	0	1	0
1	156	6	95	67
0	25	143	28	140

Table III shows that the logistic regression classifier of the testing set correctly classifies a total of 156 in class 1 and correctly classifies a total of 143 in class 0. Moreover, the k -NN algorithm of the testing set correctly classifies a total of 95 in class 1 and correctly classifies a total of 140 in class 0.

Table IV shows several comparisons of the performance evaluation results of two classification methods: the logistic regression classifier and the k -NN classifier. First, the comparison of the evaluation results on the training set for the two classifier methods is almost the same value approximation, the values of several metrics, such as

sensitivity, F1-score and AUC-ROC in the k -NN classifier are greater than the values of the metrics in the logistic regression classifiers, then for the accuracy of the two classifiers the value is the same, in other words, the evaluation results on the training set for the logistic regression classifier method is comparable to the k -NN classifier method. Next, the comparison of the evaluation results on the testing set for the two classification methods as a whole show that the value of all evaluation metrics such as sensitivity, specificity, precision, accuracy, F1-score and AUC-ROC in the logistic regression classifier method is greater than the k -NN classifier method. In other words, the evaluation results on the testing set for the logistic regression classifier method outperformed the k -NN classifier method. Overall, the performance of the logistic regression classifier method outperformed both the training set and the testing set as shown in Fig. 2.

TABLE IV
COMPARISON OF LOGISTIC REGRESSION AND K -NN EVALUATION RESULTS (%)

Technique	Evaluation	Training	Testing
Logistic Regression	Sensitivity/Recall	87.1	85.1
	Specificity	92.3	96.3
	Precision	91.7	96
	Accuracy	90	91
	F1-Score	89.4	90
	AUC-ROC	89.7	90.7
k -NN	Sensitivity/Recall	100	83.3
	Specificity	81.1	58.6
	Precision	83.8	67.6
	Accuracy	90	71
	F1-Score	91	75
	AUC-ROC	90.5	71

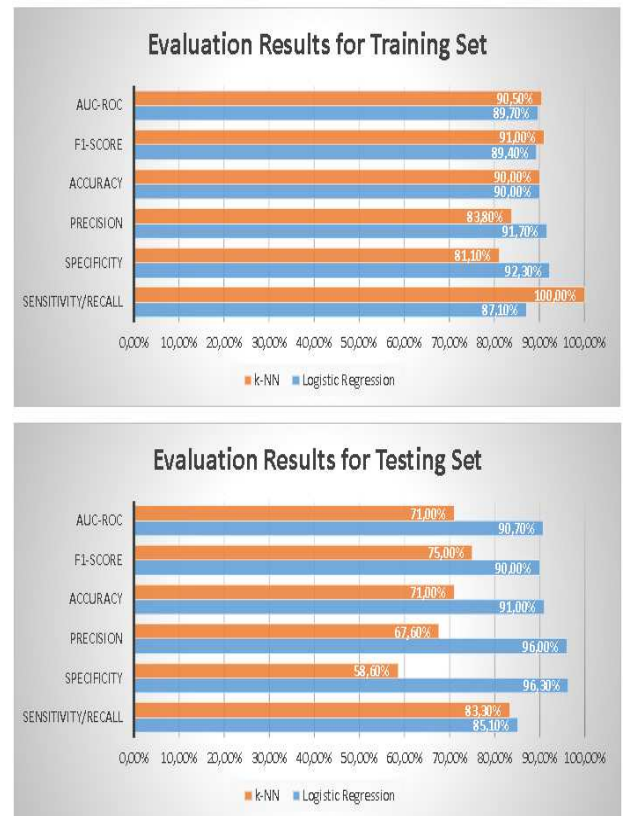


Fig. 2 Comparison of Evaluation Results for Logistic Regression and k -NN

IV. CONCLUSIONS

Predicting which customers will click on ads on websites and social media platforms is very important. The company's big goal is to target the right audience to advertise its products on websites and social media platforms. The prediction was implemented by using Logistic Regression and k -Nearest Neighbors classifiers. Results showed that the Logistic Regression model outperformed k -Nearest Neighbors model. Significant results were obtained from k -Nearest Neighbors, too, with slight differences in training sets between the models themselves, depending on the evaluation metrics. Based on the results of this study, it is recommended that further studies need to be carried out in the case of predicting customer ad clicks, such as using other types of machine learning algorithm classification techniques in order to obtain a classification method with the best performance.

REFERENCES

- [1] G. Shrivastava, V. Nagar, and S. K. Gill, "The Effects of Advertising on Consumer Buying Behavior with Special Reference to FMCG Industry," *AU-HIU Int. Multidiscip. J.*, vol. 2, no. 1, pp. 1–8, 2022.
- [2] A. Goldfarb, "What is Different About Online Advertising?," *Rev. Ind. Organ.*, vol. 44, no. 2, 2014, doi: 10.1007/s11515-013-9399-3.
- [3] R. R. Garrett, J. Yang, Q. Zhang, and S. D. Young, "An online advertising intervention to increase adherence to stay-at-home-orders during the COVID-19 pandemic: An efficacy trial monitoring individual-level mobility data," *Int. J. Appl. Earth Obs. Geoinf.*, vol. 108, 2022, doi: 10.1016/j.jag.2022.102752.
- [4] S. Gu, B. Ślusarczyk, S. Hajizada, I. Kovalyova, and A. Sakhbieva, "Impact of the covid-19 pandemic on online consumer purchasing behavior," *J. Theor. Appl. Electron. Commer. Res.*, vol. 16, no. 6, 2021, doi: 10.3390/jtaer16060125.
- [5] K. Varnali, "Online behavioral advertising: An integrative review," *Journal of Marketing Communications*, vol. 27, no. 1, 2021, doi: 10.1080/13527266.2019.1630664.
- [6] E. F. Fowler, M. M. Franz, G. J. Martin, Z. Peskowitz, and T. N. Ridout, "Political Advertising Online and Offline," *Am. Polit. Sci. Rev.*, vol. 115, no. 1, 2021, doi: 10.1017/S0003055420000696.
- [7] G. Brajnik and S. Gabrielli, "A review of online advertising effects on the user experience," *Int. J. Hum. Comput. Interact.*, vol. 26, no. 10, 2010, doi: 10.1080/10447318.2010.502100.
- [8] M. R. Farooqi and M. F. Ahmad, "The effectiveness of online advertising on consumers' mind - An empirical study," *Int. J. Eng. Technol.*, vol. 7, no. 2, 2018, doi: 10.14419/ijet.v7i2.11.11006.
- [9] S. Guha, B. Cheng, and P. Francis, "Challenges in measuring online advertising systems," 2010, doi: 10.1145/1879141.1879152.
- [10] Y. Yang and P. Zhai, "Click-through rate prediction in online advertising: A literature review," *Inf. Process. Manag.*, vol. 59, no. 2, 2022, doi: 10.1016/j.ipm.2021.102853.
- [11] P. Amlathe, "Standard Machine Learning Techniques in Audio Beehive Monitoring: Classification of Audio Samples with Logistic Regression, K-Nearest Neighbor, Random Forest and Support Vector Machine," *ProQuest Diss. Theses*, 2018.
- [12] F. Itoo, Meenakshi, and S. Singh, "Comparison and analysis of logistic regression, Naïve Bayes and KNN machine learning algorithms for credit card fraud detection," *Int. J. Inf. Technol.*, vol. 13, no. 4, 2021, doi: 10.1007/s41870-020-00430-y.
- [13] A. Kononova, W. Kim, E. Joo, and K. Lynch, "Click, click, ad: the proportion of relevant (vs. irrelevant) ads matters when advertising within paginated online content," *Int. J. Advert.*, vol. 39, no. 7, 2020, doi: 10.1080/02650487.2020.1732114.
- [14] M. Richardson, E. Dominowska, and R. Ragno, "Predicting clicks: Estimating the click-through rate for new ads," 2007, doi: 10.1145/1242572.1242643.
- [15] H. Cheng and E. Cantú-Paz, "Personalized click prediction in sponsored search," 2010, doi: 10.1145/1718487.1718531.
- [16] A. Broder *et al.*, "To swing or not to swing: Learning when (not) to advertise," 2008, doi: 10.1145/1458082.1458216.
- [17] Q. Guo and E. Agichtein, "Ready to buy or just browsing? Detecting web searcher goals from interaction data," 2010, doi: 10.1145/1835449.1835473.
- [18] D. Chakrabarti, D. Agarwal, and V. Josifovski, "Contextual advertising by combining relevance with click feedback," 2008, doi: 10.1145/1367497.1367554.
- [19] A. DeMaris, "A Tutorial in Logistic Regression," *J. Marriage Fam.*, vol. 57, no. 4, 1995, doi: 10.2307/353415.
- [20] S. Dreiseitl and L. Ohno-Machado, "Logistic regression and artificial neural network classification models: A methodology review," *J. Biomed. Inform.*, vol. 35, no. 5–6, 2002, doi: 10.1016/S1532-0464(03)00034-0.
- [21] M. R. Romadhon and F. Kurniawan, "A Comparison of Naive Bayes Methods, Logistic Regression and KNN for Predicting Healing of Covid-19 Patients in Indonesia," 2021, doi: 10.1109/EICoNCIT50028.2021.9431845.
- [22] A. E. Minarno, W. A. Kusuma, and H. Wibowo, "Performance Comparison Activity Recognition using Logistic Regression and Support Vector Machine," 2020, doi: 10.1109/ICoIAS49312.2020.9081858.
- [23] M. P. LaValley, "Logistic regression," *Circulation*, vol. 117, no. 18, pp. 2395–2399, 2008, doi: 10.1161/CIRCULATIONAHA.106.682658.
- [24] W. Cheng and E. Hüllermeier, "Combining instance-based learning and logistic regression for multilabel classification," in *Machine Learning*, 2009, vol. 76, no. 2–3, doi: 10.1007/s10994-009-5127-5.
- [25] J. Tolles and W. J. Meurer, "Logistic regression: Relating patient characteristics to outcomes," *JAMA - Journal of the American Medical Association*, vol. 316, no. 5, 2016, doi: 10.1001/jama.2016.7653.
- [26] R. Murtirawat, S. Panchal, V. K. Singh, and Y. Panchal, "Breast Cancer Detection Using K-Nearest Neighbors, Logistic Regression and Ensemble Learning," 2020, doi: 10.1109/ICESC48915.2020.9155783.
- [27] K. Shah, H. Patel, D. Sanghvi, and M. Shah, "A Comparative Analysis of Logistic Regression, Random Forest and KNN Models for the Text Classification," *Augment. Hum. Res.*, vol. 5, no. 1, 2020, doi: 10.1007/s41133-020-00032-0.
- [28] H. A. A. Rahman, Y. B. Wah, H. He, and A. Bulgiba, "Comparisons of ADABOOST, KNN, SVM and logistic regression in classification of imbalanced dataset," in *Communications in Computer and Information Science*, 2015, vol. 545, doi: 10.1007/978-981-287-936-3_6.
- [29] A. Uyar and F. Gürgeç, "Arrhythmia classification using serial fusion of support vector machines and logistic regression," 2007, doi: 10.1109/IDAACS.2007.4488483.
- [30] L. Xiong and Y. Yao, "Study on an adaptive thermal comfort model with K-nearest-neighbors (KNN) algorithm," *Build. Environ.*, vol. 202, 2021, doi: 10.1016/j.buildenv.2021.108026.
- [31] P. Verlinde and G. Cholet, "Comparing Decision Fusion Paradigms Using k-NN based Classifiers, Decision Trees and Logistic Regression in A Multi-modal Identity Verification Application," *4th Int. Conf. Audio-Video-based Biometric Pers. Authentication*, 1999.
- [32] G. A. Sandag, N. E. Tedry, and S. Lolong, "Classification of Lower Back Pain Using K-Nearest Neighbor Algorithm," 2019, doi: 10.1109/CITSM.2018.8674361.
- [33] L. M. Zouhal and T. Denoeux, "An evidence-theoretic k-NN rule with parameter optimization," *IEEE Trans. Syst. Man Cybern. Part C Appl. Rev.*, vol. 28, no. 2, 1998, doi: 10.1109/5326.669565.
- [34] M. Ali, L. T. Jung, A. H. Abdel-Aty, M. Y. Abubakar, M. Elhoseny, and I. Ali, "Semantic-k-NN algorithm: An enhanced version of traditional k-NN algorithm," *Expert Syst. Appl.*, vol. 151, 2020, doi: 10.1016/j.eswa.2020.113374.
- [35] K. Huang, S. Li, X. Kang, and L. Fang, "Spectral-Spatial Hyperspectral Image Classification Based on KNN," *Sens. Imaging*, vol. 17, no. 1, 2016, doi: 10.1007/s11220-015-0126-z.
- [36] B. Campillo-Gimenez, W. Jouini, S. Bayat, and M. Cuggia, "Improving Case-Based Reasoning Systems by Combining K-Nearest Neighbour Algorithm with Logistic Regression in the Prediction of Patients' Registration on the Renal Transplant Waiting List," *PLoS One*, vol. 8, no. 9, 2013, doi: 10.1371/journal.pone.0071991.
- [37] W. Shang, H. Huang, H. Zhu, Y. Lin, Z. Wang, and Y. Qu, "An improved kNN algorithm - Fuzzy kNN," in *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 2005, vol. 3801 LNAI, doi: 10.1007/11596448_109.
- [38] G. G. Enas and S. C. Choi, "Choice of the smoothing parameter and efficiency of k-nearest neighbor classification," *Comput. Math. with Appl.*, vol. 12, no. 2 PART A, 1986, doi: 10.1016/0898-1221(86)90076-3.

- [39] M. V. Subha and S. T. Nambi, "Classification of stock index movement using k-nearest neighbours (k-NN) algorithm," *WSEAS Trans. Inf. Sci. Appl.*, vol. 9, no. 9, 2012.
- [40] M. Saberioon, P. Čisáň, L. Labbé, P. Souček, P. Pelissier, and T. Kerneis, "Comparative performance analysis of support vector machine, random forest, logistic regression and k-nearest neighbours in rainbow trout (*Oncorhynchus mykiss*) classification using image-based features," *Sensors (Switzerland)*, vol. 18, no. 4, 2018, doi: 10.3390/s18041027.
- [41] M. E. Fischer *et al.*, "An epidemiologic study of the association between free recall dichotic digits test performance and vascular health," *J. Am. Acad. Audiol.*, vol. 30, no. 4, 2019, doi: 10.3766/jaaa.17079.
- [42] D. Chicco, N. Tötsch, and G. Jurman, "The matthews correlation coefficient (Mcc) is more reliable than balanced accuracy, bookmaker informedness, and markedness in two-class confusion matrix evaluation," *BioData Min.*, vol. 14, 2021, doi: 10.1186/s13040-021-00244-z.
- [43] S. M. Sherwood, T. B. Smith, and R. S. W. Masters, "Decision reinvestment, pattern recall and decision making in rugby union," *Psychol. Sport Exerc.*, vol. 43, 2019, doi: 10.1016/j.psychsport.2019.03.002.